

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

PARALLEL PROCESSING SYSTEM AND DATA TRANSFERRING METHOD

Patent Number: JP5173991
Publication date: 1993-07-13
Inventor(s): OKABAYASHI ICHIRO
Applicant(s): MATSUSHITA ELECTRIC IND CO LTD
Requested Patent: ☐ JP5173991
Application Number: JP19920044399 19920302
Priority Number(s):
IPC Classification: G06F15/16; G06F13/38
EC Classification:
Equivalents: JP3389610B2

Abstract

PURPOSE:To enable a PE-to-PE communication and improve processor performance.
CONSTITUTION:The parallel processing system consists of plural PEs 1a, 1c, and 1d and a network 2 which connects the PEs mutually. The PE 1a is constituted by connecting a processor 3a, a memory 4a, and a data transfer device 5a to a common bus. The data transfer device 5 has three buffers and a data repeating device 6 has two buffers. Data from the PE1a to the PE1d are transferred in the order of the memory 4a, buffer 7a, buffer 10a, buffer 8c, buffer 11e, buffer 9d, and memory 4d as shown by a dotted line. Namely, the PE1c repeats the data by utilizing the buffer 8c. Consequently, neither memory writing nor reading is performed at the repeating PE at the time of an optional PE-to- PE communication, so the overhead at the repeating PE is reduced to improve the transfer performance. Further, the data transfer device performs no bus access, so the bus width is widened and the performance of the processor is improved.

Data supplied from the esp@cenet database - I2

AN

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)特許出願公開番号

特開平5-173991

(43)公開日 平成5年(1993)7月13日

(51)Int.Cl. ⁴	識別記号	庁内整理番号	FI	技術表示箇所
G 0 6 F 15/16	4 0 0 N	9190-51		
	3 2 0 V	8840-51		
13/38	3 4 0 C	9072-5B		

審査請求 未請求 請求項の数10(全 18 頁)

(21)出願番号 特願平4-44399

(22)出願日 平成4年(1992)3月2日

(31)優先権主張番号 特願平3-54529

(32)優先日 平3(1991)3月19日

(33)優先権主張国 日本(JP)

(71)出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1000番地

(72)発明者 岡林 一郎

大阪府門真市大字門真1006番地 松下電器

産業株式会社内

(74)代理人 弁理士 小園治 明 (外2名)

(54)【発明の名称】 並列処理システムとデータ転送方法

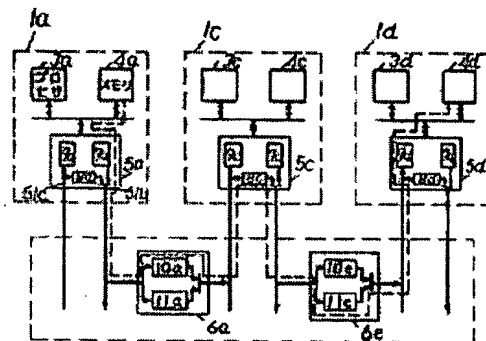
(57)【要約】

【目的】 本発明は並列処理システムに関するものであり、P E間通信及びプロセッサ性能の向上を実現する。

【構成】 並列処理システムは複数のP E 1a・1c・1dと、P E間を相互に接続するネットワーク2で構成される。P E 1aはプロセッサ3a、メモリ4a、データ転送装置5aを共通のバスに接続した構成である。データ転送装置5aは、3つのバッファを、データ中継装置6は2つのバッファを有する。P E 1aからP E 1dへは、点線で示した様に、メモリ4a、バッファ7a、バッファ10a、バッファ8a、バッファ11e、バッファ8d、メモリ4dとデータを転送する。即ち、P E 1aでバッファ8aを利用してデータを中継する。

【効果】 任意P E間通信時、中継P Eでのメモリライト/リードがないので、ここでのオーバーヘッドが軽減され、転送性能が向上する。また、ここでデータ転送装置がバスアクセスをしないので、バス幅も広がりプロセッサの性能も向上する。

1 P E Process Element
2 ネットワーク
3 プロセッサ
4 メモリ
5 データ転送装置 data 전송장치
6 データ中継装置 중계장치
7~11 バッファ



2 Network

【特許請求の範囲】

【請求項 1】 プロセッサと、メモリと、第1、第2、第3の3つのバッファを有するデータ転送装置を具備する複数のプロセッサエレメントと、前記複数のプロセッサエレメント間で直接または1つ以上のプロセッサエレメントを中継することで間接的にデータ転送が可能なように接続する

2 ネットワークを備え、データ転送の際に、送り手となるプロセッサエレメントでは、データ転送装置がメモリまたはプロセッサからデータを前記第1のバッファに取り込んだ後、前記ネットワークへ送出し、受け手となるプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを受取り、前記第2のバッファに取り込んだ後、メモリまたはプロセッサへ格納し、中継するプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを前記第3のバッファに取り込んだ後、再び前記ネットワークへ送出することを特徴とする並列処理システム。

【請求項 2】 第1、第2、第3の3つの入出力ポートと、前記第1の入出力ポートと第2の入出力ポートの間に接続された第1のバッファと、前記第1の入出力ポートと第3の入出力ポートの間に接続された第2のバッファと、前記第2の入出力ポートと前記第3の入出力ポートの間に接続された第3のバッファと、前記第1の入出力ポートから取り込んだデータにタグを付加するタグ生成部と、外部から前記第1の入出力ポートへデータを入力する際のアドレスを生成し、かつ入力回数をカウントするメモリアドレス生成部と、前記第1の入出力ポートから外部へデータを送出する際にデータ出力回数をカウントするカウンタと、前記第1の入出力ポートから出力するデータの一部分をアドレスとして取り出し、データ出力の際にはこれを、データ入力の際には前記アドレス生成部のアドレスを選択して外部に出力する第1のセクタと、前記第2の入出力ポートから外部へデータを送出する際に、外部アドレスを生成する第2の中継アドレス生成部と、前記第3の入出力ポートへ外部からデータを入力する際に、外部アドレスを生成する第2の中継アドレス生成部と、第2の入出力ポートから外部へデータを送出する際、前記第1のバッファと前記第3のバッファ出力の一部分を選択する第2のセクタと、前記第1のバッファからデータを送出する場合には前記第1の中継アドレス生成部の出力を、前記第3のバッファからデータを送出する場合には前記第3のバッファ出力の他の一部分を選択する第3のセクタとを具備することを特徴とするデータ転送装置。

【請求項 3】 第1、第2、第3の3つの入出力ポートと、前記第1の入出力ポートと第2の入出力ポートの間に接続された第1のバッファと、前記第1の入出力ポートと第3の入出力ポートの間に接続された第2のバッファと、前記第2の入出力ポートと前記第3の入出力ポートの間に接続された第3のバッファと、前記第1の入出力

ポートから取り込んだデータにタグを付加するタグ生成部と、外部から前記第1の入出力ポートへデータを入力する際のアドレスを生成し、かつ入力回数をカウントするメモリアドレス生成部と、前記第1の入出力ポートから外部へデータを送出する際にデータ出力回数をカウントするカウンタと、前記第1の入出力ポートから出力するデータの一部分をアドレスとして取り出し、データ出力の際にはこれを、データ入力の際には前記アドレス生成部のアドレスを選択して外部に出力する第1のセクタと、第2の入出力ポートから外部へデータを送出する際、前記第1のバッファ出力の一部分と前記第3のバッファ出力の一部分を選択する第2のセクタと、前記第1のバッファ出力の他の一部分と前記第3のバッファ出力の他の一部分を選択する第3のセクタとを具備することを特徴とするデータ転送装置。

【請求項 4】 第1、第2、第3の3つの入出力ポートと、第1の入出力ポートと第2の入出力ポートの間に接続された第1のバッファと、前記第1の入出力ポートと第3の入出力ポートの間に接続された第2のバッファと、前記第2の入出力ポートと前記第3の入出力ポートの間に接続された第3のバッファと、前記第1の入出力ポートから取り込んだデータにタグを付加するタグ生成部と、外部から前記第1の入出力ポートへデータを入力する際のアドレスを生成し、かつ入力回数をカウントするメモリアドレス生成部と、前記第1の入出力ポートから外部へデータを送出する際にデータ出力回数をカウントするカウンタと、前記第1の入出力ポートから出力するデータの一部分をアドレスとして取り出し、データ出力の際にはこれを、データ入力の際には前記アドレス生成部のアドレスを選択して外部に出力する第1のセクタと、前記第2の入出力ポートから外部へデータを送出する際に、外部アドレスを生成する第1の中継アドレス生成部と、前記第3の入出力ポートから外部へデータを送出する際に、外部アドレスを生成する第2の中継アドレス生成部と、第2の入出力ポートから外部へデータを送出する際、前記第1のバッファと前記第3のバッファ出力の一部分を選択する第2のセクタと、前記第1のバッファからデータを送出する場合には前記第1の中継アドレス生成部の出力を、前記第3のバッファからデータを送出する場合には前記第3のバッファ出力の他の一部分を選択する第3のセクタと、第3の入出力ポートから外部へデータを送出する際、前記第2のバッファと前記第3のバッファ出力の一部分を選択する第4のセクタと、前記第2のバッファからデータを送出する場合には前記第2の中継アドレス生成部の出力を、前記第3のバッファからデータを送出する場合には前記第3のバッファ出力の他の一部分を選択する第5のセクタを具備することを特徴とするデータ転送装置。

【請求項 5】 少なくとも3つの入出力ポートと、第1の入出力ポートと第2の入出力ポートの間に接続された第

1のバッファと、前記第1の入出力ポートと第3の入出力ポートの間に接続された第2のバッファと、前記第2の入出力ポートと前記第3の入出力ポートの間に接続された第3のバッファと、前記第1の入出力ポートから取り込んだデータにタグを付加するタグ生成部と、外部から前記第1の入出力ポートへデータを入力する際のアドレスを生成し、かつ入力回数をカウントするメモリアドレス生成部と、前記第1の入出力ポートから外部へデータを送出する際にデータ出力回数をカウントするカウンタと、前記第1の入出力ポートから出力するデータの一部をアドレスとして取り出し、データ出力の際にはこれを、データ入力の際には前記アドレス生成部のアドレスを選択して外部に出力する第1のセレクタと、第2の入出力ポートから外部へデータを送出する場合、前記第1のバッファ出力の一部と前記第3のバッファ出力の一部を選択する第2のセレクタと、前記第1のバッファ出力の他の一部と前記第3のバッファ出力の他の一部を選択する第3のセレクタと、第3の入出力ポートから外部へデータを送出する場合、前記第2のバッファ出力の一部と前記第3のバッファ出力の一部を選択する第4のセレクタと、前記第2のバッファ出力の他の一部と前記第3のバッファ出力の他の一部を選択する第5のセレクタとを具備することを特徴とするデータ転送装置。

【請求項 6】 請求項 2又は4記載のタグ生成部は、データの種別を判定する制御情報部と、2回目以降に中継する中継部のアドレスを中継部に示す複数の中継アドレス部と、最終的なデータの格納アドレスを示すメモリアドレス部を有する情報タグとして付加することを特徴とするデータ転送装置。

【請求項 7】 請求項 3又は5記載のタグ生成部は、データの種別を判定する制御情報部と、中継する中継部のアドレスを中継部に示す複数の中継アドレス部と、最終的なデータの格納アドレスを示すメモリアドレス部を有する情報タグとして付加することを特徴とするデータ転送装置。

【請求項 8】 第1,第2の2つの入出力ポートと、 N を2以上の整数として N 個のバッファと、 N 入力1出力の出力セレクタと、 $N-1$ 個の2入力1出力の入力セレクタを有し、前記出力セレクタの入力に前記 N 個のバッファ出力を接続し、出力を前記第1の入出力ポートに接続し、前記第2の入出力ポートを第1のバッファの入力及び前記 $N-1$ 個の入力セレクタの入力的一端に接続し、第2から第 N までのバッファ入力に前記入力セレクタの出力を接続し、 L を1以上 $N-1$ 以下の整数として第 L のバッファ出力を第 L の入力セレクタの他端に接続することを特徴とするデータ中継装置。

【請求項 9】 プロセッサと、メモリと、請求項 2,3のいずれかに記載のデータ転送装置を共通のバスに接続した構成の複数のプロセッサエレメントと、各格子点に請求項8記載のデータ中継装置を配してここを経由することで

ある2つのプロセッサエレメント間通信を行い、また N を2以上の整数として、 $N-1$ 回プロセッサエレメントを中継することで任意のプロセッサエレメント間通信が可能となるネットワークを備え、データ転送の際に、送り手となるプロセッサエレメントでは、データ転送装置がメモリまたはプロセッサからデータを第1のバッファに取り込んだ後前記ネットワークへ送出し、中継するプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを第3のバッファに取り込んだ後、再び前記ネットワークへ送出し、受け手となるプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを第2のバッファに取り込んだ後、メモリまたはプロセッサへ格納し、また、 L を1以上 $N-1$ の整数として、第 L 回目のプロセッサエレメントから第 $L+1$ 回目のプロセッサエレメントへのデータ転送時に、前記データ中継装置の第 L のバッファを経由することを特徴とする並列処理システム。

【請求項 10】 N を2以上の整数として、少なくとも2つのポートを有する N 個のプロセッサエレメントと、 K, L を1以上 N 以下の整数として、 $N \times N$ 個の格子点を有し、各格子点を (K, L) とし、この格子点に少なくとも2つのポートを有するバッファ (K, L) を配したネットワークを備え、前記第 K のプロセッサエレメントの一端をバッファ (K, L) の一端に共通に接続し、また前記バッファ (K, L) の他端を L が共通になるように接続し、この共通接続線を前記プロセッサエレメントに接続するか、または外部ポートとする構成を基本単位として含む並列処理システムにおいて、前記第 K のプロセッサエレメントからデータを送出する際に、バッファ (K, L) から順次データを送出するデータ転送方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、計算機分野でその将来性が期待されている並列処理システムに係わり、特にプロセッサエレメント間通信に関する。

【0002】

【従来の技術】 一般に並列処理システムは、計算処理を行なうプロセッサエレメント（以下PE）と、PE間のデータ転送を行なうネットワークで構成される。

【0003】 以下図面を参照しながら、従来の並列処理システムにおけるPE間通信の一例について説明する。図10は従来の第1の並列処理システムの構成図、図11は従来のデータ転送装置の構成図、図12は従来のデータ中継装置の構成図を示すものである。これらは、電子情報通信学会・集積回路研究会ICD89-152に詳しく述べられている。ここでは、PE、ネットワークの一部を示す。また、データの流れを単方向に限定して説明する。

【0004】 まず、従来の並列処理システム（図10）について説明する。基本的には、PE1a, 1b, 1dとPE間を相互に接続するネットワーク2で構成される。PE1

は全て同一の構成であり、P E1aを例にとれば、プロセッサ3a、メモリ4a、データ転送装置5aを共通のバスに接続した構成である。また、データ転送装置5aは2つのバッファ7a、9aを有する。また、ネットワーク2内部にデータ中継装置6a、6eを設ける。データ中継装置6a・6eはそれぞれバッファ10a、10eを有する。ネットワーク2は、任意P E間通信が第3のP Eを1回経由することで可能なもの

(P E間距離が2)である。以上の従来の並列処理システムにおいて、P E1aからP E1dへのデータの流は、メモリ4a、バッファ7a、バッファ10a、バッファ9e、メモリ4e、バッファ7e、バッファ10e、バッファ9d、メモリ4dとなる。これを図10中に点線で示す。

【0005】続いて、従来のデータ転送装置(図11)について説明する。入力ポート17aはメモリ4に、出力ポート17bはネットワーク2に接続される。入力ポート17aから17bへのデータの流は次の様になる。メモリアドレス生成部12aよりセクタ18a経由で、アドレス50aを出力してメモリリードを行ないデータ51aを入力ポート17aからバッファ7へ取り込む。次に、中継アドレス生成部15aからアドレス50bを出力して、出力ポート17bからデータ51bを出力する。

【0006】また、出力ポート17cから17aへのデータの流は次の様になる。中継アドレス生成部15bからアドレス50cを出力して、出力ポート17cからデータ51cを入力し、バッファ9に取り込む。次に、出力ポート17aから、メモリアドレス生成部12bよりセクタ18a経由でアドレス50aを、バッファ9よりデータ51aを出力してメモリライトを行なう。なお、制御部16a・16bはバッファ状態52a・52bを監視する。

【0007】さらに続いて、従来のデータ中継装置(図12)について説明する。データ51bはバッファ10に格納される。制御部31aはバッファ10のリード/ライトを制御する。デコーダ30a・30bはアドレス50b・50cを監視し、自分がアクセスされた際に、トライステートバッファ32a・32cをイネーブルとして、バッファ状態52a・52bを外部へ通過させる。ここでのバッファ状態とは、書き込み側はバッファフル、読みだし側はバッファエンプティに関するものである。

【0008】次に、図13は従来のデータ転送方法を示す図である。これは、ネットワーク2が完全クロスバ網の例である。P Eからのデータ送出順序をデータ中継装置6a-6p内に示す。即ち、最初のステップでP E1aはデータ中継装置6aに、1bは6e、1cは6i、1dは6mに一斉にデータを送出する。次のステップでは、P E1aはデータ中継装置6bに、1bは6f、1cは6j、1dは6nに一斉にデータを送出する。以下同様で、端まで送出し終わると最初に戻る。最初のステップ終了後、P E1aがデータ中継装置6aよりデータを受信する。

【0009】最後に、図14は従来の第2の並列処理システムの構成図である。これは、電子情報通信学会・コ

ンピュータシステム 研究会CPSY89-11に詳しく述べられている。ここで、P U(プロセッシングユニット)は図14(a)に示す様にメッシュ状に接続されている。各P Uは図14(b)に示す様にC P U71、ローカルメモリ72、周辺L S I73を共通のバスに接続した構成である。また、4つのポート75a-dを有し、2ポートR A Mであるコネクションメモリ74a-bを介して他P Uと通信を行なう。

【0010】

【発明が解決しようとする課題】しかしながら上記の様な第1の並列処理システムでは、中継するP Eがデータを一旦メモリにライトしてから再度リードするので、ここでのオーバーヘッドが大きいという問題点を有していた。またメモリアクセスを行なうので、バスネックが発生しプロセッサ性能も低下する。

【0011】また、上記の様なデータ転送方式では、最初のステップ終了後P E1aのみがデータ中継装置6a、6e、6i、6mに接続されているため受信可能で、このバスだけに負荷が集中し、システム全体の転送性能は低下する。

【0012】また、上記の様な第2の並列処理システムでは、全P Uが同期して隣接するP Uと通信する場合は非常に高速であるが、距離の遠いP Uとの通信は遅い。任意のP E間距離は $N \times N$ のシステムで最大 N 、平均 $N/2$ である。また個々のP Uで通信要求がランダムに発生する場合の対応、あるいは他ネットワークへの拡張の2点でも不利である。

【0013】本発明は上記問題点に鑑み、高いプロセッサ性能、高速なP E間通信を実現し、かつ柔軟な並列処理システムを提供することを目的とする。

【0014】

【課題を解決するための手段】上記問題点を解決するために、本発明の並列処理システムは、プロセッサと、メモリと、第1、第2、第3の3つのバッファを有するデータ転送装置を具備する複数のプロセッサエレメントと、前記複数のプロセッサエレメント間で直接または1つ以上のプロセッサエレメントを中継することで間接的にデータ転送が可能なように接続するネットワークを備えたものであり、データ転送の際に、送り手となるプロセッサエレメントでは、データ転送装置がメモリまたはプロセッサからデータを第1のバッファに取り込んだ後から前記ネットワークへ送出し、受け手となるプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを第2のバッファに取り込んだ後、メモリまたはプロセッサへ格納し、中継するプロセッサエレメントでは、データ転送装置が前記ネットワークから前記データを第3のバッファに取り込んだ後、再び前記ネットワークへ送出するものである。

【0015】また、本発明のデータ転送方式は、少なくとも2つのポートを有する N 個(N は2以上の整数)のプロセッサエレメントと、 $N \times N$ 個の格子点を有し、各格

予点を(K, L) (K, Lは1以上N以下の整数)とし、ここに少なくとも2つのポートを有するバッファ(K, L)を配したネットワークを備え、前記第Kのプロセサエレメントの一端をバッファ(K, L)の一端に共通に接続し、またバッファ(K, L)の他端をLが共通になるように接続し、この共通接続線を前記プロセサエレメントに接続するか、または外部ポートとする構成を基本単位として含む並列処理システムにおいて、前記第Kのプロセサエレメントからデータを送出する際に、バッファ(K, K)から順次データを送出するものである。

【0016】

【作用】本発明の並列処理システムでは、上記した構成によって、PE間で通信を行なう場合に、送り手となるプロセサエレメントでは、データ転送装置がメモリまたはプロセサからデータを第1のバッファに取り込んでからネットワークへ送出する。また、受け手となるプロセサエレメントでは、データ転送装置がネットワークからデータを第2のバッファに取り込んだ後、メモリまたはプロセサへ格納する。さらに、中継するプロセサエレメントでは、データ転送装置がネットワークからデータを第3のバッファに取り込んだ後、再びネットワークへ送出する。従って、中継PEでのメモリアイト/リードがないので、ここでのオーバーヘッドが軽減され、転送性能が向上する。また、データ転送装置がバスアクセスをしないので、バス幅も広がりプロセサの性能も向上する。

【0017】また、本発明のデータ転送方式では、 $N \times N$ 個(Nは2以上の整数)の格子点を有するクロスバ網に対して、第Kのプロセサエレメントは、格子点(K, K)から順次データを送出するので、受信の際特定のバスにトラフィックが集中することはなく、転送性能が向上する。

【0018】

【実施例】以下本発明の実施例について、図面を参照しながら説明する。

【0019】これらの実施例では説明及び図面の簡化のため、データの流れを一方向に限定する。また、ネットワークは、任意PE間通信が第3のPEを1回経由することで可能なもの(PE間距離が2)である。

【0020】本発明の第1の実施例における並列処理システムの全体構成図(図6)を説明する。これはPE数4の並列処理システムである。図6の構成によりPE間距離2を実現する。図1はこの抜粋と考えてよく、PE1a、データ中継装置6a、PE1c、データ中継装置6c、PE1dと接続されることがわかる。

【0021】図1は本発明の第1の実施例における並列処理システムの構成図である。基本的には、PE1a・1c・1dとPE間を相互に接続するネットワーク2で構成される。制御線を含めた全体的な接続は、図7を用いて後ほど説明する。

【0022】PEは全て同一の構成であり、PE1aを例

にとれば、プロセサ3a、メモリ4a、データ転送装置5aを共通のバスに接続した構成である。また、データ転送装置5aは3つのバッファ7a・8a・9aを有する。また、ネットワーク2内部にデータ中継装置6a・6cを設ける。データ中継装置6aはバッファ10a・11aを、6cはバッファ10c・11cを有する。以上の並列処理システムにおいて、PE1aからPE1dへのデータの流れは、メモリ4a、バッファ7a、バッファ10a、バッファ8c、バッファ11c、バッファ9d、メモリ4dとなる。これを図1中に点線で示す。

【0023】ここで、データの中継するPE1cにおいて、データ転送装置5cはデータ中継装置6aから受信したデータを内部のバッファ8c経由でデータ中継装置6cに送出する。ここで、メモリ4cのアクセスを伴わないので、転送速度の向上及びバス幅の拡大の両者が同時に可能となる。

【0024】以下、本実施例を実現するための要素技術であるデータ転送装置、データ中継装置等について順次説明する。

【0025】まず、第1のデータ転送装置について説明する。図2は、本発明の第1の実施例におけるデータ転送装置の構成図である。入出力ポート17aはメモリ4c、17b、17cはネットワーク2のバッファ中継装置6cに接続される。

【0026】入出力ポート17aから17bへのデータの流れは次のようになる。メモリアドレス生成部12aよりセクタ18a経由で、アドレス50を出力してメモリアイトを行なう。データ51を入出力ポート17aから取り込み、タグ生成部13の出力をタグとしてデータに付加した後、バッファ7へ取り込む。次に、中継アドレス生成部15aからセクタ18c経由でアドレス50b、バッファ7からセクタ18b経由でデータ51bを入出力ポート17bより出力する。ここで、メモリアドレス生成部12aはリード回数のカウントも行なう。タグについては、図4を用いて後ほど説明する。

【0027】また、入出力ポート17cから17aへのデータの流れは次のようになる。中継アドレス生成部15bからアドレス50cを出力して、入出力ポート17cからデータ51cを入力し、バッファ9に取り込む。次に、入出力ポート17aから、バッファ9の出力の一部をセクタ18a経由でアドレス50に、バッファ9の他の一部をデータ51に出力してメモリアイトを行なう。ここで、カウンタ14はライト回数のカウントを行なう。

【0028】最後に、入出力ポート17cから17bへのデータの流れは次のようになる。中継アドレス生成部15bからアドレス50cを出力して、入出力ポート17cからデータ51cを入力し、タグ変換部131でタグ部分を変換したのちバッファ8に取り込む。次に、アドレス50bとしてバッファ8の出力の一部をセクタ18c経由で、データ51bとしてバッファ8の出力の他の一部をセクタ18b経由でそれぞれ入出力ポート17bより出力する。即ち中継アドレス

として、メモリからネットワークへの転送時は中継アドレス生成部15a出力を、ネットワークからネットワークへの転送時はデータの一部を用いる。またメモリアドレスとして、リード時はメモリアドレス生成部18a出力を、ライト時はデータの一部を用いる。

【0029】なお、制御部16a、16bはデータ中継装置のバッファ状態52a、52bを監視する。ここでは、入出力ポート17b、17cを単方向としたが、これは双方向でもよい。図15はこれを示したもので、第1の実施例におけるデータ転送装置を双方向にした場合の構成図である。この場合はバッファ7、8、9、内部線を双方向化した上で、セクタ18a、18dを入出力ポート17c側に設ける。タグ交換部131は1つであるので、データの流れが入出力ポート17c→入出力ポート17bの場合はバッファ8入力時、入出力ポート17b→入出力ポート17cの場合はバッファ9出力時にタグの交換をすることになる。

【0030】次に、第2のデータ転送装置について説明する。図3は本発明の第2の実施例におけるデータ転送装置の構成図である。基本的には図2と同様であるが、入出力ポート17aから17bへのデータの流に若干の相違点があるのでこれについて説明する。メモリアドレス生成部12aよりセクタ18a経由で、アドレス50を出力してメモリリードを行なう。データ51を入出力ポート17aから取り込み、タグ及び中継アドレス生成部130の出力をタグとしてデータに付加した後、バッファ7へ取り込む。次に、アドレス50bとしてバッファ7の出力の一部をセクタ18c経由で、データ51bとしてバッファ7の他の一部をセクタ18b経由でそれぞれ入出力ポート17bより出力する。即ち中継アドレスとして、メモリからネットワークへの転送時、ネットワークからネットワークへの転送時にデータの一部を用いる。図2ではデータ中継装置15aで生成したアドレスを、図3ではタグ及び中継アドレス生成部130で生成することになる。

【0031】ここでは、入出力ポート17b、17cを単方向としたが、これは双方向でもよい。図16はこれを示したもので、第2の実施例におけるデータ転送装置を双方向にした場合の構成図である。この場合はバッファ7、8、9、内部線を双方向化した上で、セクタ18a、18dを入出力ポート17c側に設ける。タグ交換部131は1つであるので、データの流れが入出力ポート17c→入出力ポート17bの場合はバッファ8入力時、入出力ポート17b→入出力ポート17cの場合はバッファ9出力時にタグの交換をすることになる。

【0032】次に、タグ生成部13で生成するタグについて、本実施例におけるデータ形式の構成図である図4を用いて説明する。データはデータ部24とタグ部で構成される。タグは、データの種別（属性・形式等）を示す制御情報部20、現在の転送が第何ステップであるかを示す回数部21、2回目以降に中継する中継部のアドレスを中継順に示す複数の中継アドレス部22a、22b、最終的なデ

ータの格納アドレスを示すメモリアドレス部23を備える。図4(a)が送出、中継P/E間、図4(b)が中継、受信P/E間のデータの形式である。図1の例では図4(a)がP/E1a-1c間、図4(b)がP/E1c-1d間のデータの形式である。

【0033】図1を用いてタグとP/E間通信の関係について説明する。データ転送装置は図3の構成とする。データ転送装置5aのタグ及び中継アドレス生成部130で、制御情報部20に制御情報を、回数部21には回数を、中継アドレス部22aにはデータ中継装置6aのアドレスを、中継アドレス部22bにはデータ中継装置6eのアドレスを、メモリアドレス部23にはメモリ4dのアドレスをそれぞれ示したタグを生成・付加する。これが図4(a)である。例として、制御情報部20はデータ長を示す3、回数部21は一回目を示す1、中継アドレス部22aはデータ中継装置6aが手前にあるので0、中継アドレス部22bも0を付加する。データ転送装置5aは、中継アドレス部22aの“0”を送出してデータ中継装置6aへデータを送出する。ここで制御情報部20の3はデータ長が2の3乗つまり8ワードである旨等を含む。次にデータ転送装置5cのタグ交換部131ではデータ中継装置6aからデータを取り込み、回数部21を2回目を示す2にし、中継アドレス部22aを削除して図4(b)に示す形式に変換後、中継アドレス部22bの“0”を送出してデータ中継装置6eへデータを送出する。最後に、データ転送装置5dではデータ中継装置6eからデータを取り込んだ後、メモリアドレス部23を出力してメモリ4dへデータをライトする。データ部が8ワードであるので、メモリアドレス部23に示すアドレスより順次8ワードライトすることになる。

【0034】ここで、データ取り込みに関して、候補となるデータ中継装置が複数ある場合は、順次スキャンしてデータの準備されたものから取り込みめばよい。これら一連の動作で、データ転送装置はネットワークから受けたデータを再度ネットワークに送出するか、メモリにライトするか決める必要がある。これは例えば、回数部21を見て判断すればよい。また、別の方法としてデータ中継装置にP/E間距離に相当する複数のバッファを準備し、ある特定のバッファからのデータはメモリ、その他はネットワークと決めれば回数部21は不要となる。これは、図5にて後ほど説明する。

【0035】また、データ転送装置が図2の構成の場合には、図4の中継アドレス部22aは存在せず、代わりに中継アドレス生成部15aが一回目の中継アドレスを生成することになる。なお、タグの構成としては回数部21の代わりに、複数の中継アドレス部の最後に終り符号を付ける形式も可能である。

【0036】さらに別の方法として、タグ交換をせず図4(a)の形式のままで、1回目は22aを、2回目は22bを中継アドレスとして送出する方法もある。この場合は図2等におけるタグ交換部131は不要であるが、セク

タ18cで1、2回目で選択するビット位置を変える必要が生じる。また少しではあるが、タグのビット数が大きくなる。

【0037】データ毎にアドレス情報が付加された以上の構成により、送り手及び受け手のPEが複数で、かつデータ数が複数で、また流れる順序がランダムな場合でも、複雑な制御なしで、確実に転送が実現できる。

【0038】さて次に、データ中継装置について説明する。図5は本発明の実施例におけるデータ中継装置の構成図である。データ中継装置には2つのモードを有する。

【0039】第1のモードは、1本のバッファとして動作するもので、入力セクタ35はバッファ10出力を、出力セクタ34はバッファ11出力を選択する。バッファ10及び11を連続した1本のバッファとして使用する。動作は従来例(図12)と同様である。即ち、制御部31aはバッファ10、11のリード/ライトを制御する。デコーダ30a・30bはアドレス50b・50cを監視し、自分がアクセスされた際に、トライステートバッファ32a・32bをイネーブルとして、バッファ状態52a・52bを外部へ通過させる。この時、制御部31b、トライステートバッファ32c・32dはディセーブルである。

【0040】第2のモードは、2本のバッファを並列に動作させるもので、入力セクタ35は入出力ポート36aを、出力セクタ34は必要に応じてバッファ10または11出力を選択する。バッファ10及び11を独立な2本のバッファとして使用する。

【0041】制御部31aはバッファ10のリード/ライトを、制御部31bはバッファ11のリード/ライトを制御する。デコーダ30a・30bはアドレス50b・50cを監視し、自分がアクセスされた際に、トライステートバッファ32a-32dをイネーブルとして、バッファ状態52a-52dを外部へ通過させる。ここでのバッファ状態とは、書き込み側はバッファフル、読み出し側はバッファエンプティに関するものであり、バッファ状態52a・52bがバッファ10に、バッファ状態52c・52dがバッファ11にそれぞれ対応する。

【0042】第2のモードでは、データ中継装置6にPE間距離2に相当する2本のバッファが存在することになる。データ中継装置6a・6eでは、バッファ10a・10eが1回目、バッファ11a・11eが2回目のデータを格納する。従って、データは前記した様に図1の点線の流れとなる。

【0043】データ転送装置5は、データ中継装置の2つのバッファ状態を監視して、送受可能な方とデータのやりとりを行なう。ここで、データ転送装置5cがバッファ10aから取り込んだ時はバッファ8cに、バッファ11aから取り込んだ時はバッファ9cにそれぞれ格納する制御を行なうことで、図4のタグの回数部21は不要となる。

【0044】第3のPEを介するPE間通信時、本デ

ータ中継装置の第1のモードを用いるデータ中継装置を用いた場合はデッドロックが発生するが、本データ中継装置の第2のモードを用いることで、1回目と2回目のデータが独立に扱えるので、デッドロックが回避できるが、この事情について説明する。

【0045】図17は本発明のデータ中継装置(第1のモード)を用いた場合の転送の様子を示す図、図18は本発明のデータ中継装置(第2のモード)を用いた場合の転送の様子を示す図である。ここでは、PE1b、1a、1e間でデータが流れる場合について考える。また簡単化のためバッファ7、8、9は1段、図17でバッファ10d、10a、10fは2段、図18でバッファ10d、10a、10f、11d、11a、11fは1段とする。

【0046】送信、中継、受信PEとデータの関係は次の様になる。

送信PE1b→中継PE1a→受信PE1e : データc1, c2, c3, c4

送信PE1a→中継PE1c→受信PE1b : データb1, b2, b3, b4

送信PE1c→中継PE1b→受信PE1a : データa1, a2, a3, a4

データ中継装置(第1のモード)を用いた場合では、図17に示す様な状態に陥った場合にデッドロックとなる。ここで例えばPE1bはデータc4またはa2を送出したいが、データ中継装置6dのバッファ10dがフルであるので送れない。データ中継装置6dはデータc3を吐き出したいがc3が入るべきバッファ8aがフルであるので転送できない。バッファ8aに空きが生じるためにはバッファ10aに空きが生じる必要があるが、データb3が入るべきバッファ8cがフルであるのでバッファ10aは空かない。バッファ8cが空くためにはバッファ10fが空く必要があるが、このためにはバッファ8bが空く必要がある。そのためにはバッファ10dが空く必要があり、結局どのバッファも空くことはない、つまりデッドロックとなる。この様に複数のPEで閉じたループを構成する場合にデッドロックが発生する可能性が高い。

【0047】データ中継装置(第2のモード)を用いた場合では、図17に相当する状態が図18(a)である。例えばデータ中継装置6dについてみればバッファ10dに1回目の転送途中のデータc3、バッファ11dに2回目の転送途中のデータa1が格納される。

【0048】次のサイクルではデータa1がバッファ9a、データc1がバッファ9c、データb1がバッファ9bに転送される。この状態を示したのが図18(b)である。バッファ9a、9c、9bのデータはメモリにライトされるのでこれらのバッファにはすぐ空きが生じる。こうなると例えばPE1bはデータa2をバッファ11d経由でバッファ9aに送れる。バッファ8bに空きが生じるのでバッファ10fのデータa3がバッファ8bに転送可能となる。以下同様に順次データが流れデッドロックは生じない。

【0049】以上により、ランダムな通信要求が発生した場合でも確実に動作できる。また1回目と2回目のデータの優先度付けを適切に行なうことで転送性能も向上する。また直接P E間で転送する場合は、第1のモードで大きなバッファリングが可能となる。

【0050】データ転送装置とデータ中継装置間の制御線を含めた信号線の接続の様子を図7に示す。データ中継装置6aと6b、6aと6cの信号線が共通に接続される。中継アドレスによりアクティブなデータ中継装置が選択されてデータ・アドレスが受け渡される。またバッファ状態は選択されたデータ中継装置のみが出力し、他のデータ中継装置はハイインピーダンスである。

【0051】次にネットワークでのアドレス・データの形式を、本発明の実施例におけるアドレス・データ構成図である図8を用いて説明する。ここでは、図8(a)、(b)2つの例を示す。一般に並列処理システムではメモリとネットワークでのバス幅が異なり、ネットワーク側が狭くなる。そのため、ネットワークとのインターフェースでデータ幅の変換が必要となる。

【0052】図8(a)では、データ転送装置1aからデータ中継装置6a間のデコーダ30aへアドレス50bが入力される。データ51b(アドレス以外)は、データ転送装置1aの出力ラッチ40に格納された後、セクタ41で分解されて、データ中継装置6aへ入力される。図8(b)では、アドレス50b及びデータ51bはデータ転送装置1aの出力ラッチ40に格納された後、セクタ41で分解されて、データ中継装置6aへ入力され、アドレス50bはデコーダ30aへ、データ51bは内部へ入力される。図8(b)は、データにアドレス情報が含まれ、データ中継装置は入力データを常に監視し、自分のアドレスに対応するものが出現した場合にデータを取り込むデータフロー的な制御となる。図8(a)に比べて、制御ロジックは複雑になり、転送量も多いが、データ転送装置とデータ中継装置間の配線数は少なくなる。

【0053】最後に、本発明の実施例におけるデータ転送方法について図9、図19で説明する。図9は本発明の実施例におけるデータ転送方法を示す図、図19は同実施例における時間と転送レートの関連図である。図9は、ネットワーク2が完全クロスバ網の例である。P Eからのデータ送出順序をデータ中継装置6a-6p内に示す。即ち、最初のステップでP E1aはデータ中継装置6aに、1bは6f、1cは6k、1dは6pに一番にデータを送出する。次のステップでは、P E1aはデータ中継装置6bに、1bは6g、1cは6l、1dは6mに一番にデータを送出する。以下同様で、端まで送出し終わると最初に戻る。これにより最初のステップ終了後、全てのP Eで受信が可能となる。即ち、P E1aはデータ中継装置6a、1bは6f、1cは6k、1dは6pよりそれぞれデータを受信できる。従って特定のネットワークの負荷が偏らないので、システム全体の転送効率が向上する。これを図19に示す。従来は、

時間1で1つのチャネルでのみ送受が行なわれるので転送レートは1である。時間2で転送レート2、3で3、4で4と増えて行きその後減少し、時間7で終了する。これを点線で示す。本実施例では実線で示した様に時間1-4で全チャネルが動作つまり転送レートが4であり、時間4で転送は終了する。

【0054】なお、本例は既に説明した図6の様に部分的にクロスバを有するシステムに適用可能である。また、最初に述べた様に、ここではデータの流れを一方に限定したが、バッファの双方向化、セクタなど一部の回路を2つ持つことで、二方向の流れにも容易に対応できる。

【0055】また、P E間距離が2の並列処理システムについて述べたが、タグの中継アドレス、データ中継装置のバッファ本数をNとすることで、P E間距離がNの並列処理システムに拡張可能である。またこれらを組み合わせて、各種の形態のネットワークを実現することが可能となる。

【0056】

【発明の効果】以上述べてきた様に、本発明の並列処理システムでは、任意P E間通信時、中継P Eでのメモリライト/リードがないので、ここでのオーバーヘッドが軽減され、転送性能が向上する。また、データ転送装置がバスアクセスをしないので、バス幅も広がりプロセサの性能も向上する。

【0057】また本発明のデータ転送方式では、ネットワークの負荷が分散する、つまり受信側のP Eが均等に動作できるので、システム全体の転送性能が向上する。

【0058】さらに、本発明のデータ転送装置、データ中継装置を用いることで、各種のネットワークを有する並列処理システムが構成できる。

【0059】単体プロセサの計算機性能及び半導体技術の限界が見えてきた現在、並列処理システムへの期待は非常に大きく、本発明は極めて有用なものである。

【図面の簡単な説明】

【図1】本発明の第1の実施例における並列処理システムの構成図

【図2】本発明の第1の実施例におけるデータ転送装置の構成図

【図3】本発明の第2の実施例におけるデータ転送装置の構成図

【図4】本発明の実施例におけるデータ形式の構成図

【図5】本発明の実施例におけるデータ中継装置の構成図

【図6】本発明の第1の実施例における並列処理システムの全体構成図

【図7】同実施例における接続詳細図

【図8】本発明の実施例におけるアドレス・データ構成図

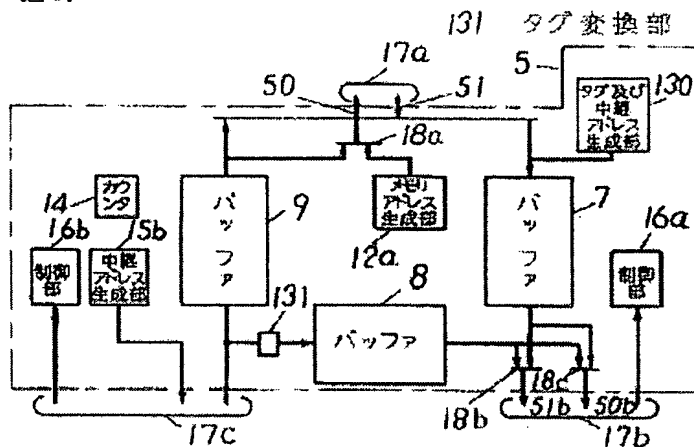
【図9】本発明の実施例におけるデータ転送方法を示す

図

- 【図 10】従来の第 1 の並列処理システム の構成図
 【図 11】従来のデータ転送装置の構成図
 【図 12】従来のデータ中継装置の構成図を示す図
 【図 13】従来のデータ転送方法を示す図
 【図 14】従来の第 2 の並列処理システム の構成図
 【図 15】第 1 の実施例におけるデータ転送装置を双方向にした場合の構成図
 【図 16】第 1 の実施例におけるデータ転送装置を双方向にした場合の構成図
 【図 17】データ中継装置（第 1 のモード）を用いた場合の転送の様子を示す図
 【図 18】データ中継装置（第 2 のモード）を用いた場合の転送の様子を示す図
 【図 19】本発明の実施例におけるデータ転送方法における時間と転送レートの関連図
 【符号の説明】
 1 PE（プロセサエレメント）
 2 ネットワーク
 3 プロセサ
 4 メモリ
 5 データ転送装置
 6 データ中継装置
 7-11 バッファ
 12 メモリアドレス生成部
 13 タグ生成部
 14 カウンタ
 15 中継アドレス生成部

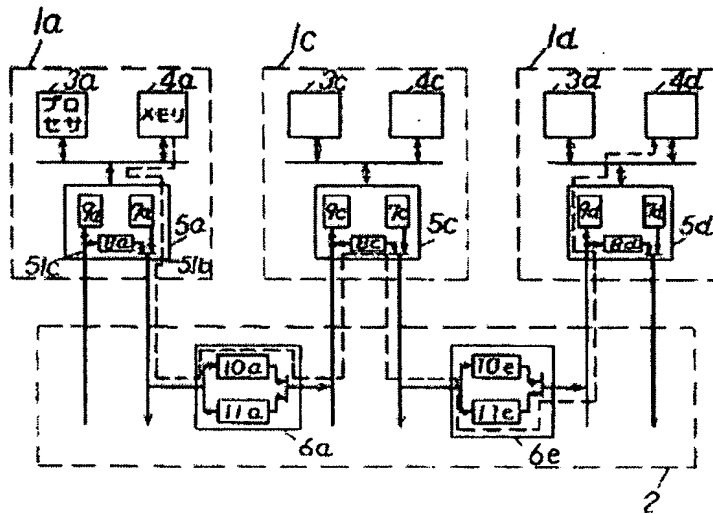
- 16 制御部
 17 入出力ポート
 18 セレクタ
 20 制御情報部
 21 回数部
 22 中継アドレス部
 23 メモリアドレス部
 24 データ部
 30 デコーダ
 31 バッファ制御部
 32 トライステートバッファ
 34 出力セレクタ
 35 入力セレクタ
 36 入出力ポート
 40 出力ラッチ
 41 セレクタ
 50 アドレス
 51 データ
 52 バッファ状態
 70 PU（プロセッシングユニット）
 71 CPU
 72 ローカルメモリ
 73 周辺LSI
 74 コネクションメモリ
 75 ポート
 130 タグ及び中継アドレス生成部
 131 タグ交換部

【図 3】

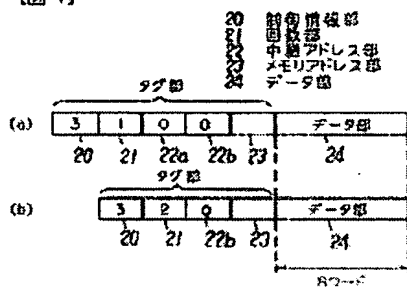


【図1】

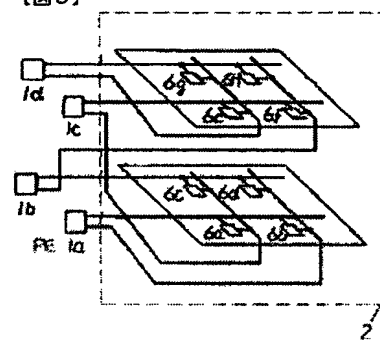
- 1 P E
- 2 ネットワーク
- 3 プロセサ
- 4 メモリ
- 5 データ転送装置
- 6 データ中継装置
- 7~11 バッファ



【図4】

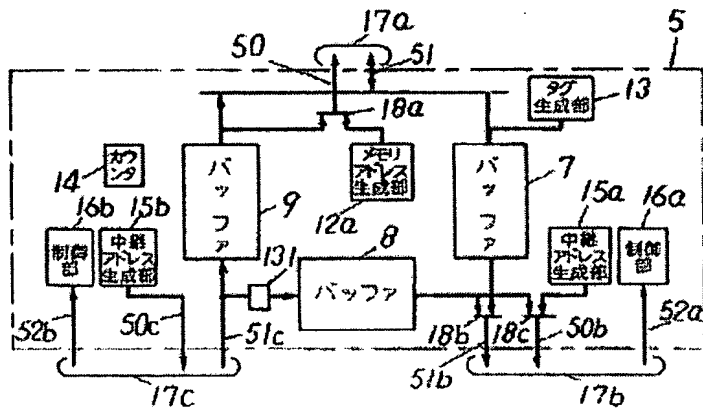


【図6】



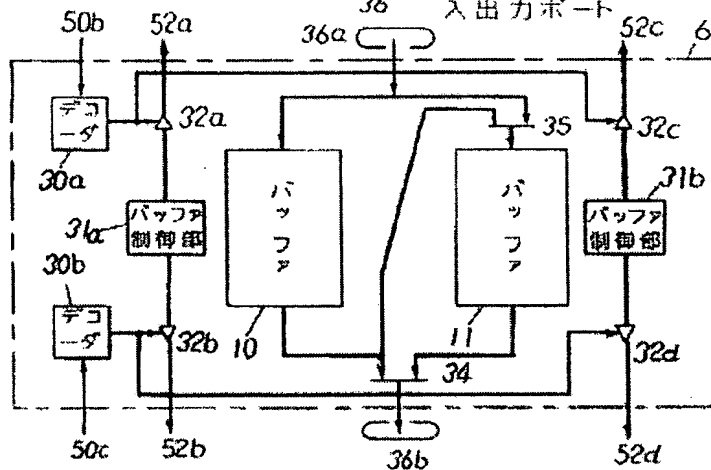
【図 2】

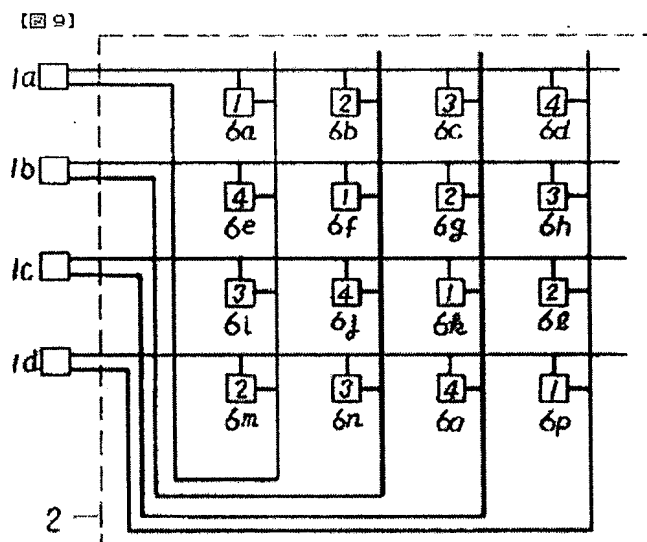
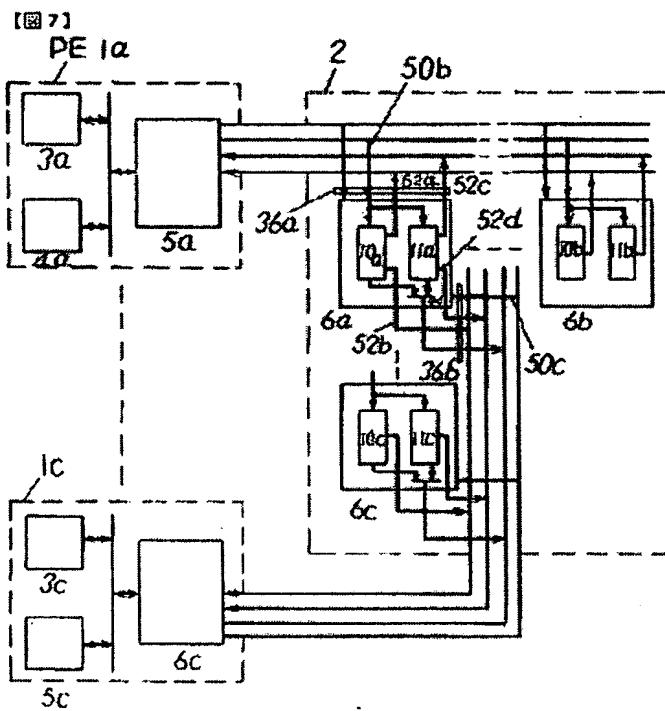
- 17 入出力ポート
- 18 セレクタ
- 50 アドレス
- 51 データ
- 52 バッファ状態



【図 5】

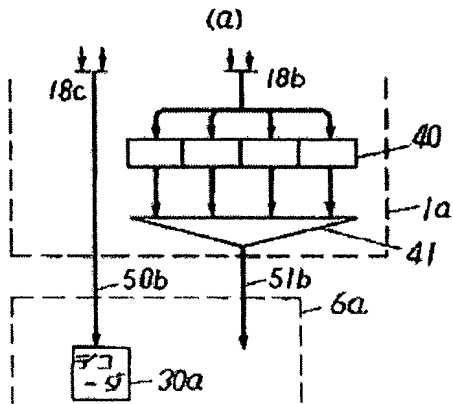
- 32 トライステートバッファ
- 34 出力セレクタ
- 35 入力セレクタ
- 36 入出力ポート



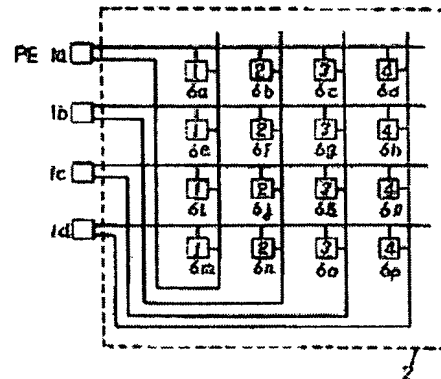


【図 9】

40 出力ラッチ
41 セレクタ



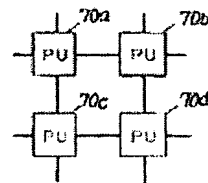
【図 13】



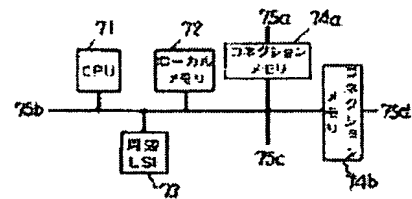
【図 14】

75 ポート

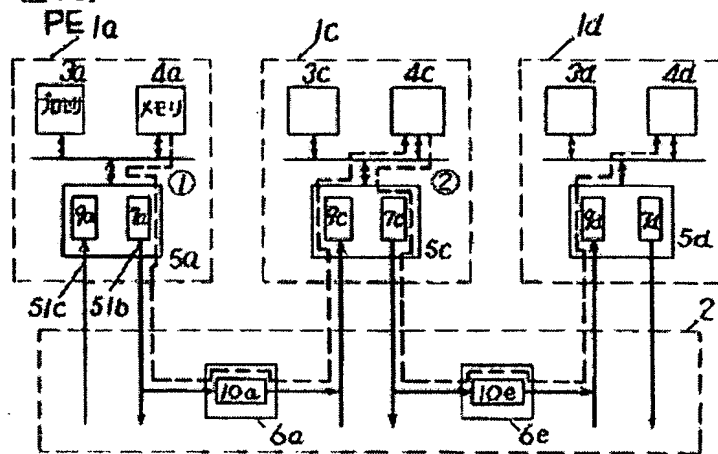
(a)



(b)

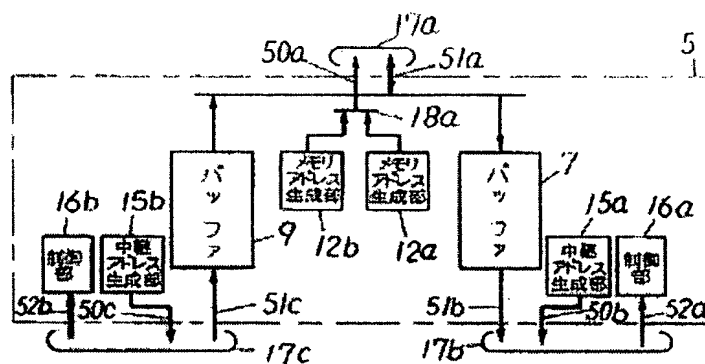


【図10】

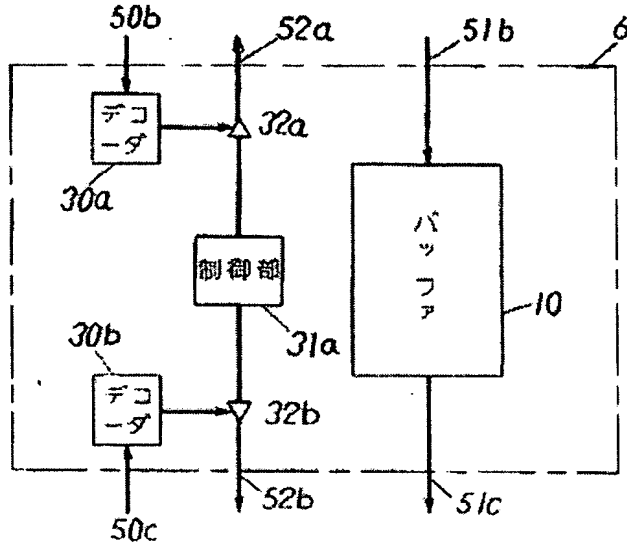


【図11】

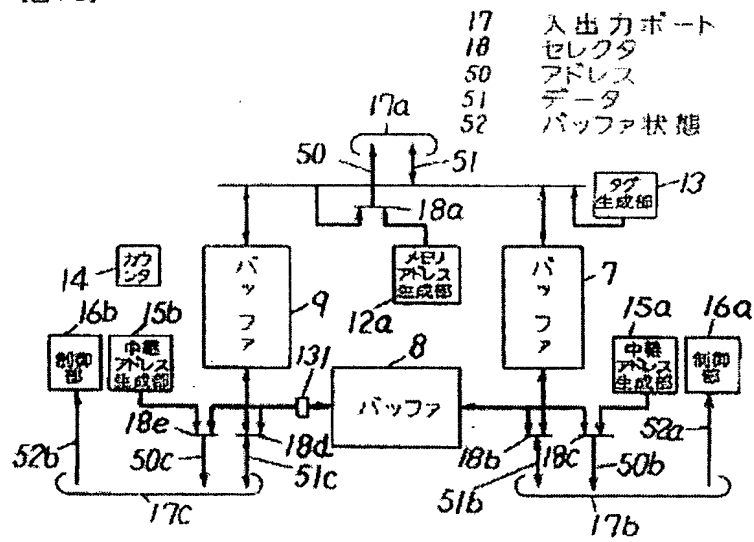
17a, 17b, 17c 入出力ポート

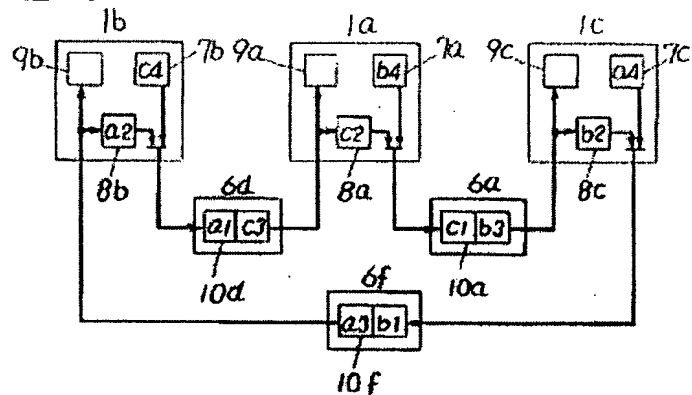


【図12】

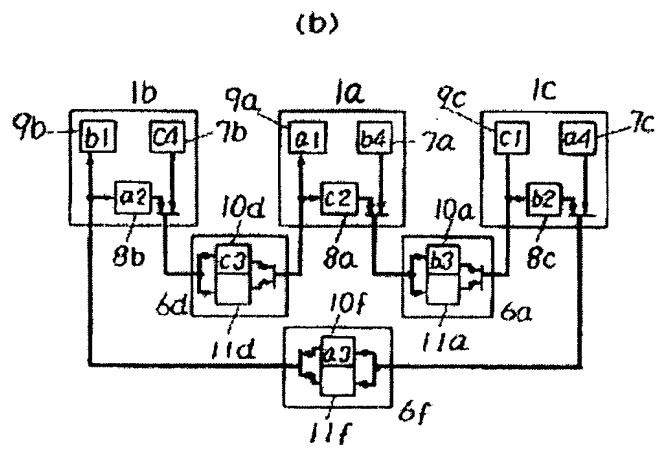
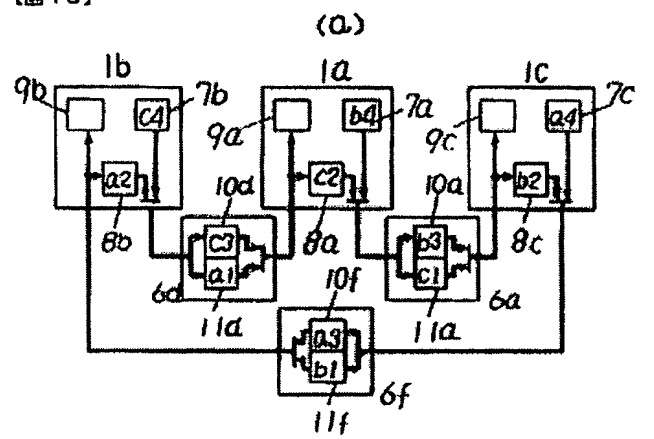


【図15】





【図18】



〔図19〕
転送レート

